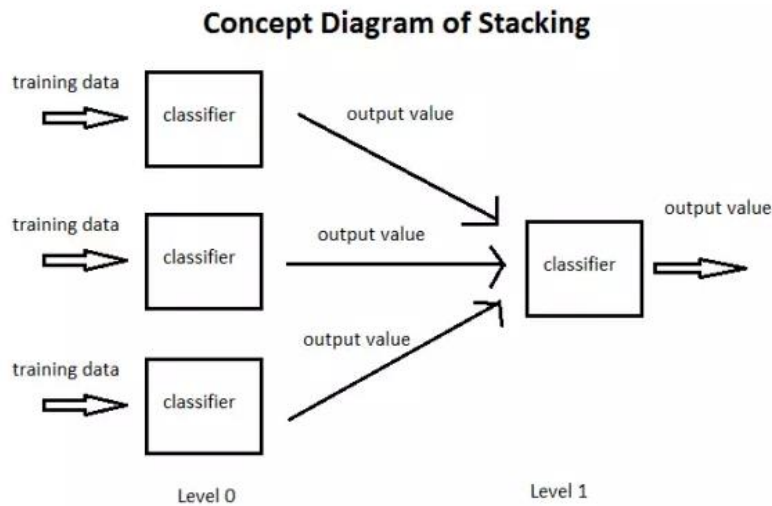


以前真的没用过这个技术

1) 什么是 stacking

简单来说 stacking 就是当用初始训练数据学习出若干个基学习器后，将这几个学习器的预测结果作为新的训练集，来学习一个新的学习器。



将个体学习器结合在一起的时候使用的方法叫做结合策略。对于分类问题，我们可以使用投票法来选择输出最多的类。对于回归问题，我们可以将分类器输出的结果求平均值。

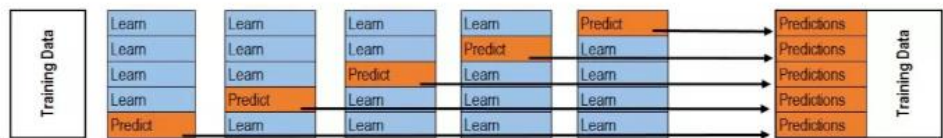
上面说的投票法和平均法都是很有效的结合策略，还有一种结合策略是使用另外一个机器学习算法来将个体机器学习器的结果结合在一起，这个方法就是Stacking。

在stacking方法中，我们把个体学习器叫做初级学习器，用于结合的学习器叫做次级学习器或元学习器（meta-learner），次级学习器用于训练的数据叫做次级训练集。次级训练集是在训练集上用初级学习器得到的。

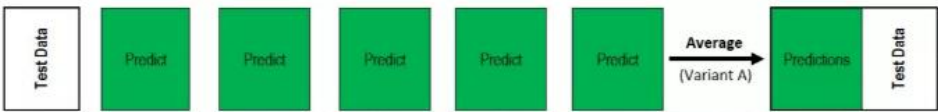
Stacking本质上就是这么直接的思路，但是直接这样有时对于如果训练集和测试集分布不那么一致的情况下是有一点问题的，其问题在于用初始模型训练的标签再利用真实标签进行再训练，毫无疑问会导致一定的模型过拟合训练集，这样或许模型在测试集上的泛化能力或者说效果会有一定的下降，因此现在的问题变成了如何降低再训练的过拟合性，这里我们一般有两种方法。

- 1. 次级模型尽量选择简单的线性模型
- 2. 利用K折交叉验证

K-折交叉验证：训练：



预测：



2) Stacking融合(回归):

```
1 from sklearn import linear_model
2
3 def Stacking_method(train_reg1,train_reg2,train_reg3,y_train_true,test_pre1,test_pre2,test_pre3,model_L2= linear_model.LinearRegression()):
4     model_L2.fit(pd.concat([pd.Series(train_reg1),pd.Series(train_reg2),pd.Series(train_reg3)],axis=1).values,y_train_true)
5     Stacking_result = model_L2.predict(pd.concat([pd.Series(test_pre1),pd.Series(test_pre2),pd.Series(test_pre3)],axis=1).values)
6     return Stacking_result
```

```
1 ## 生成一些简单的样本数据, test_pre1 代表第1个模型的预测值
2 train_reg1 = [3.2, 8.2, 9.1, 5.2]
3 train_reg2 = [2.9, 8.1, 9.0, 4.9]
4 train_reg3 = [3.1, 7.9, 9.2, 5.0]
5 # y_test_true 代表第模型的真正值
6 y_train_true = [3, 8, 9, 5]
7
8 test_pre1 = [1.2, 3.2, 2.1, 6.2]
9 test_pre2 = [0.9, 3.1, 2.0, 5.9]
10 test_pre3 = [1.1, 2.9, 2.2, 6.0]
11
12 # y_test_true 代表第模型的真正值
13 y_test_true = [1, 3, 2, 6]
```

```
1 model_L2= linear_model.LinearRegression()
2 Stacking_pre = Stacking_method(train_reg1,train_reg2,train_reg3,y_train_true,
3                               test_pre1,test_pre2,test_pre3,model_L2)
4 print('Stacking_pre MAE:',metrics.mean_absolute_error(y_test_true, Stacking_pre))
```

Stacking_pre MAE: 0.04213483146067476

可以发现模型结果相对于之前有进一步的提升,这是我们需要关注的一点是,对于第二层Stacking的模型不宜选取的过于复杂,这样会导致模型在训练集上过拟合,从而使得在测试集上并不能达到很好的效果。

Voting投票机制:

Voting即投票机制,分为软投票和硬投票两种,其原理采用少数服从多数的思想。

```
1 硬投票:对多个模型直接进行投票,不区分模型结果的相对重要度,最终投票数最多的类为最终被预测的类。
2 '''
3 iris = datasets.load_iris()
4
5 x=iris.data
6 y=iris.target
7 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
8
9 clf1 = XGBClassifier(learning_rate=0.1, n_estimators=150, max_depth=3, min_child_weight=2, subsample=0.7,
10                    colsample_bytree=0.6, objective='binary:logistic')
11 clf2 = RandomForestClassifier(n_estimators=50, max_depth=1, min_samples_split=4,
12                             min_samples_leaf=63, oob_score=True)
13 clf3 = SVC(C=0.1)
14
```

```

15 # 硬投票
16 eclf = VotingClassifier(estimators=[('xgb', clf1), ('rf', clf2), ('svc',
clf3)], voting='hard')
17 for clf, label in zip([clf1, clf2, clf3, eclf], ['XGBBoosting', 'Random
Forest', 'SVM', 'Ensemble']):
18     scores = cross_val_score(clf, x, y, cv=5, scoring='accuracy')
19     print("Accuracy: %0.2f (+/- %0.2f) [%s]" % (scores.mean(),
scores.std(), label))

```

```

1 '''
2 软投票：和硬投票原理相同，增加了设置权重的功能，可以为不同模型设置不同权重，进而
区别模型不同的重要度。
3 '''
4 x=iris.data
5 y=iris.target
6 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
7
8 clf1 = XGBClassifier(learning_rate=0.1, n_estimators=150, max_depth=3, mi
n_child_weight=2, subsample=0.8,
9     colsample_bytree=0.8, objective='binary:logistic')
10 clf2 = RandomForestClassifier(n_estimators=50, max_depth=1, min_samples_
split=4,
11     min_samples_leaf=63, oob_score=True)
12 clf3 = SVC(C=0.1, probability=True)
13
14 # 软投票
15 eclf = VotingClassifier(estimators=[('xgb', clf1), ('rf', clf2), ('svc',
clf3)], voting='soft', weights=[2, 1, 1])
16 clf1.fit(x_train, y_train)
17
18 for clf, label in zip([clf1, clf2, clf3, eclf], ['XGBBoosting', 'Random
Forest', 'SVM', 'Ensemble']):
19     scores = cross_val_score(clf, x, y, cv=5, scoring='accuracy')
20     print("Accuracy: %0.2f (+/- %0.2f) [%s]" % (scores.mean(),
scores.std(), label))

```

分类的Stacking\Blending融合：

stacking是一种分层模型集成框架。

以两层为例，第一层由多个基学习器组成，其输入为原始训练集，第二层的模型则是以第一层基学习器的输出作为训练集进行再训练，从而得到完整的stacking模型，stacking两层模型都使用了全部的训练数据。

```

1 '''
2 5-Fold Stacking

```

```

3  '''
4  from sklearn.ensemble import RandomForestClassifier
5  from sklearn.ensemble import ExtraTreesClassifier, GradientBoostingClassifier
6  import pandas as pd
7  #创建训练的数据集
8  data_0 = iris.data
9  data = data_0[:100,:]
10
11  target_0 = iris.target
12  target = target_0[:100]
13
14  #模型融合中使用到的各个单模型
15  clfs = [LogisticRegression(solver='lbfgs'),
16          RandomForestClassifier(n_estimators=5, n_jobs=-1, criterion='gini'),
17          ExtraTreesClassifier(n_estimators=5, n_jobs=-1, criterion='gini'),
18          ExtraTreesClassifier(n_estimators=5, n_jobs=-1, criterion='entropy'),
19          GradientBoostingClassifier(learning_rate=0.05, subsample=0.5,
20                                     max_depth=6, n_estimators=5)]
21
22  #切分一部分数据作为测试集
23  X, X_predict, y, y_predict = train_test_split(data, target, test_size=0.3, random_state=2020)
24
25  dataset_blend_train = np.zeros((X.shape[0], len(clfs)))
26  dataset_blend_test = np.zeros((X_predict.shape[0], len(clfs)))
27
28  #5折stacking
29  n_splits = 5
30  skf = StratifiedKFold(n_splits)
31  skf = skf.split(X, y)
32
33  for j, clf in enumerate(clfs):
34      #依次训练各个单模型
35      dataset_blend_test_j = np.zeros((X_predict.shape[0], 5))
36      for i, (train, test) in enumerate(skf):
37          #5-Fold交叉训练，使用第i个部分作为预测，剩余的部分来训练模型，获得其预测的输出作为第i部分的新特征。
38          X_train, y_train, X_test, y_test = X[train], y[train], X[test], y[test]
39          clf.fit(X_train, y_train)
40          y_submission = clf.predict_proba(X_test)[: , 1]
41          dataset_blend_train[test, j] = y_submission

```

```

41 dataset_blend_test_j[:, i] = clf.predict_proba(X_predict)[:, 1]
42 #对于测试集，直接用这k个模型的预测值均值作为新的特征。
43 dataset_blend_test[:, j] = dataset_blend_test_j.mean(1)
44 print("val auc Score: %f" % roc_auc_score(y_predict,
dataset_blend_test[:, j]))
45
46 clf = LogisticRegression(solver='lbfgs')
47 clf.fit(dataset_blend_train, y)
48 y_submission = clf.predict_proba(dataset_blend_test)[:, 1]
49
50 print("Val auc Score of Stacking: %f" % (roc_auc_score(y_predict, y_subm
ission)))

```

Blending，其实和Stacking是一种类似的多层模型融合的形式

其主要思路是把原始的训练集先分成两部分，比如70%的数据作为新的训练集，剩下30%的数据作为测试集。

在第一层，我们在这70%的数据上训练多个模型，然后去预测那30%数据的label，同时也预测test集的label。

在第二层，我们就直接用这30%数据在第一层预测的结果做为新特征继续训练，然后用test集第一层预测的label做特征，用第二层训练的模型做进一步预测

其优点在于：

- 1.比stacking简单（因为不用进行k次的交叉验证来获得stacker feature）
- 2.避开了一个信息泄露问题：generlizers和stacker使用了不一样的数据集

缺点在于：

- 1.使用了很少的数据（第二阶段的blender只使用training set10%的量）
- 2.blender可能会过拟合
- 3.stacking使用多次的交叉验证会比较稳健 "

可以发现 基模型 用 'KNN', 'Random Forest', 'Naive Bayes' 然后再这基础上 次级模型加一个 'LogisticRegression'，模型测试效果有着很好的提升。

经验总结¶

比赛的融合这个问题，个人的看法来说其实涉及多个层面，也是提分和提升模型鲁棒性的一种重要方法：

- 1) 结果层面的融合，这种是最常见的融合方法，其可行的融合方法也有很多，比如根据结果的得分进行加权融合，还可以做Log, exp处理等。在做结果融合的时候，有一个很重要的条件是模型结果的得分要比较近似，然后结果的差异要比较大，这样的结果融合往往有比较好的效果提升。
- 2) 特征层面的融合，这个层面其实感觉不叫融合，准确说可以叫分割，很多时候如果我们用同种模型训练，可以把特征进行切分给不同的模型，然后在后面进行模型或者结果融合有时也能产生比较好的效果。
- 3) 模型层面的融合，模型层面的融合可能就涉及模型的堆叠和设计，比如加Staking层，部分模型的结果作为特征输入等，这些就需要多实验和思考了，基于模型层面的融合最好不同模型类型要有一定的差异，用同种模型不同的参数的收益一般是比较小的。

[https://tianchi.aliyun.com/notebook-ai/detail?
spm=5176.20850271.J_3678908510.21.5ef03daaAHPwjA&postId=170929](https://tianchi.aliyun.com/notebook-ai/detail?spm=5176.20850271.J_3678908510.21.5ef03daaAHPwjA&postId=170929)